

PR #38191 完整报告

vllm-project/vllm

[Bugfix] Fix k_norm weight sharding in MiniMaxM2Attention when total_num_kv_heads < tp_size

合并时间: 2026-04-27 20:57

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/38191>

执行摘要

- 一句话: 修复 MiniMaxM2 在 KV head 数少于 TP 大小时 k_norm 权重分片错误
- 推荐动作: 该 PR 值得阅读以了解 TP 下权重分片的细节, 特别是 `weight_shard_world_size` 参数的设计, 对于其他需要自定义分片的模块有参考价值。建议在类似场景 (如分组 query attention) 中复用此模式。

功能与动机

当 KV head 数量 (`total_num_kv_heads`) 小于 Tensor Parallel 大小 (`tp_size`) 时, KV head 会在多个 rank 之间复制。这种情况下, MiniMaxText01RMSNormTP 的权重分片原先错误地使用 `tp_world` 而非 `total_num_kv_heads`, 导致运行时形状不匹配。参考 PR body: 'When the number of KV heads is less than the tensor parallel size, KV heads are replicated across ranks. In this case, MiniMaxText01RMSNormTP was incorrectly sharding k_norm weights by `tp_world` instead of `total_num_kv_heads`, causing a shape mismatch at runtime.'

实现拆解

1. 重构 MiniMaxText01RMSNormTP 权重分片逻辑: 在 `vllm/model_executor/layers/mamba/linear_attn.py` 中, 为 `__init__` 添加可选参数 `weight_shard_world_size` 和 `weight_shard_rank`, 允许调用者覆盖默认的 `tp_world` 和 `tp_rank`。 `weight_loader` 改为支持这些新参数, 并使用 `functools.partial` 在构造时绑定分片配置。同时将 `CustomOp` 注册从类属性 `name` 改为 `@CustomOp.register` 装饰器。
2. 调整 MiniMaxM2Attention 中 k_norm 的创建: 在 `vllm/model_executor/models/minimax_m2.py` 的 `MiniMaxM2Attention.__init__` 中, 根据 `total_num_kv_heads` 与 `tp_size` 的比较分情况创建 `k_norm`: 当大于等于 `tp_size` 时保持原有行为; 当小于时, 计算 `num_kv_head_replicas = tp_size // total_num_kv_heads`, 并传入 `weight_shard_world_size=total_num_kv_heads` 和 `weight_shard_rank=当前 rank // num_kv_head_replicas`, 确保权重按 head 数均匀分片。
3. 修复 `fuse_minimax_qk_norm` 编译融合: 原实现中 `k_norm.weight` 的形状错误, 导致 `torch.compile` 的融合 `pass` 无法匹配模式。修复后权重形状正确, 融合可以正常执行, 提升推理性能。
4. 验证方式: 通过 `lm_eval` 在 MiniMax-M2.5 模型上运行 `gsm8k` 任务, 修复后准确率与修复前一致, 且无运行时错误。

关键文件:

- vllm/model_executor/layers/mamba/linear_attn.py (模块 线性注意力; 类别 source; 类型 data-contract; 符号 init, weight_loader) : 核心修改文件: 重构 MiniMaxText01RMSNormTP, 引入 weight_shard_world_size 和 weight_shard_rank 参数, 修正 CustomOp 注册方式。
- vllm/model_executor/models/minimax_m2.py (模块 模型定义; 类别 source; 类型 data-contract) : 调用侧调整文件: 在 MiniMaxM2Attention 中根据 total_num_kv_heads 与 tp_size 的关系动态创建 k_norm, 传入合适的分片参数。

关键符号: MiniMaxText01RMSNormTP.init, MiniMaxText01RMSNormTP.weight_loader, MiniMaxM2Attention.init

关键源码片段

vllm/model_executor/layers/mamba/linear_attn.py

核心修改文件: 重构 MiniMaxText01RMSNormTP, 引入 weight_shard_world_size 和 weight_shard_rank 参数, 修正 CustomOp 注册方式。

```
# vllm/model_executor/layers/mamba/linear_attn.py
# MiniMaxText01RMSNormTP: 支持自定义权重分片维度的 RMSNorm 层

@CustomOp.register("minimax_text01_rmsnorm_tp")
class MiniMaxText01RMSNormTP(CustomOp):
    def __init__(
        self,
        hidden_size: int,
        eps: float = 1e-6,
        *,
        weight_shard_world_size: int | None = None, # 允许调用者指定分片维度
        weight_shard_rank: int | None = None,
    ) -> None:
        super().__init__()
        self.tp_world = get_tensor_model_parallel_world_size()
        self.tp_rank = get_tensor_model_parallel_rank()
        # 若不传入则默认为 tp_world / tp_rank, 保持原有行为
        self.weight_shard_world = weight_shard_world_size or self.tp_world
        self.weight_shard_rank = (
            self.tp_rank if weight_shard_rank is None else weight_shard_rank
        )

        # 使用整数除法 // 避免浮点结果
        self.weight = nn.Parameter(torch.ones(hidden_size // self.weight_shard_world))
        # 将分片参数绑定到 weight_loader, 避免全局变量依赖
        self.weight.weight_loader = partial(
            self.weight_loader,
            shard_world_size=self.weight_shard_world,
            shard_rank=self.weight_shard_rank,
        )
```

```

self.variance_epsilon = eps

@staticmethod
def weight_loader(
    param: nn.Parameter,
    loaded_weight: torch.Tensor,
    shard_world_size: int | None = None, # 注入的分片世界大小
    shard_rank: int | None = None,
) -> None:
    # 如果未提供则回退到全局 tp 信息
    if shard_world_size is None:
        shard_world_size = get_tensor_model_parallel_world_size()
    if shard_rank is None:
        shard_rank = get_tensor_model_parallel_rank()

    shard_size = loaded_weight.shape[0] // shard_world_size
    shard = slice(shard_rank * shard_size, (shard_rank + 1) * shard_size)
    param.data.copy_(loaded_weight[shard])

```

vllm/model_executor/models/minimax_m2.py

调用侧调整文件：在 MiniMaxM2Attention 中根据 total_num_kv_heads 与 tp_size 的关系动态创建 k_norm，传入合适的分片参数。

```

# vllm/model_executor/models/minimax_m2.py
# MiniMaxM2Attention.__init__ 中 k_norm 的创建分支

# 原有逻辑：直接用 tp_size 分片
if self.total_num_kv_heads >= tp_size:
    self.k_norm = MiniMaxText01RMSNORMTP(
        self.head_dim * self.total_num_kv_heads, eps=rms_norm_eps
    )
else:
    # 当 KV head 数少于 TP 大小时，KV head 会被复制到多个 rank 上。
    # 权重应按 KV head 数分片，而非 tp_size。
    num_kv_head_replicas = tp_size // self.total_num_kv_heads
    self.k_norm = MiniMaxText01RMSNORMTP(
        self.head_dim * self.total_num_kv_heads,
        eps=rms_norm_eps,
        weight_shard_world_size=self.total_num_kv_heads, # 按 head 数分片
        weight_shard_rank=get_tensor_model_parallel_rank() // num_kv_head_replicas,
    )

```

评论区精华

Review 中的核心讨论：

- tdoublep 问为什么移除 name = "MiniMaxText01RMSNORMTP"；作者回复正确做法是使用 @CustomOp.register 装饰器。

- ZJY0516 质疑是否真的需要 CustomOp（因为没有不同实现需要分发）；作者解释类已继承自 CustomOp，此修改只是纠正注册方式，保持正确性。
- jeejeelee 询问此变更是否影响 fuse_minimax_qk_norm 融合 pass；作者说明原代码下权重形状错误导致融合失败，修复后权重形状正确，融合可以正常工作。
- 移除 name 属性改用 @CustomOp.register (design): 同意使用装饰器方式注册 CustomOp，确保正确注册。
- 是否需要 CustomOp 类型 (design): 保持 CustomOp 注册，因为类已经继承自 CustomOp，修改确保正确注册。
- 对 fuse_minimax_qk_norm 的影响 (correctness): 修复不仅没有破坏融合，反而使融合在低 KV head 场景下正常工作。

风险与影响

- 风险：风险较低，限定在 MiniMaxM2 模型及其 RMSNormTP 层。引入的新参数有默认回退行为（若未传入则与旧行为一致），未引入向后兼容性问题。functools.partial 的使用在层序列化场景下需注意，但通常不会序列化模型层。主要风险是 TP 配置下代码路径必须正确匹配，但通过手动验证确认。缺少自动化测试覆盖，但作者已通过 lm_eval 端到端验证。
- 影响：影响范围：使用 MiniMaxM2 模型且在 Tensor Parallel 配置下，特别是当 KV head 数量少于 TP 大小（例如 num_kv_heads=8, tp_size=16）的用户得到修复，运行时不再出现形状错误。对于无 TP 或 KV head 数不小于 TP 的用户无行为变化。此外，fuse_minimax_qk_norm 编译融合的恢复可带来性能提升。影响程度为中低，仅针对特定模型和配置。
- 风险标记：权重分片逻辑变更，CustomOp 注册变更，缺少自动化测试

关联脉络

- 暂无明显关联 PR